

A Bayesian Analysis of the Multinomial Randomized Response Model Using Dirichlet Prior Distribution

Jong-Min Kim¹

Statistics, Division of Science and Mathematics, University of Minnesota, Morris, MN, 56267, USA

Tae-Young Heo²

Department of Statistics, North Carolina State University, Raleigh, NC, 27695, USA

Abstract

In this paper, we examine the problem of estimating the sensitive characteristics and behaviors in a multinomial randomized response (RR) model. We analyze this problem through a Bayesian perspective and develop a Bayesian multinomial RR model in survey study. The Bayesian inference of multinomial RR model is a new approach to RR models.

Keywords: Randomized response; Multinomial model; Dirichlet distribution

1 Introduction

The frequency of socially undesirable, embarrassing, or prohibited acts or attitudes is usually underestimated in surveys. A randomized response (RR) technique is a procedure for collecting the information on sensitive characteristics without exposing the identity of the respondent. RR technique was originally proposed by Warner (1965) as an alternative survey technique for socially undesirable or incriminating behavior questions. With the many benefits of Dirichlet prior in Bayesian framework, we propose a Bayesian multinomial approach to an extension of the binomial randomized response model suggested by Kim, Tebbs and An (2005).

¹ Jong-Min Kim, Assistant Professor of Statistics, Division of Science and Mathematics, University of Minnesota, Morris, MN, 56267, USA Email: jongmink@morris.umn.edu

² Tae-Young Heo, Ph.D. Candidate, Department of Statistics, North Carolina State University, Raleigh, NC, 27695, USA Email: heoty@hanmail.net

2 Multinomial Randomized Response Model

Using the Hopkins randomized device, Kim and Warde (2005) propose a multinomial RR model and derive estimators and their properties. We follow Kim and Warde's (2005) multinomial model set up which explicitly assumed a multinomial model for a single sensitive variable, denoted as A . Suppose that there are two different colors of balls, red and green, in the device and that each of the green balls contains a discrete number $1, 2, \dots, k$. All green balls represent a set of non-sensitive categories, $B = \{B_1, B_2, \dots, B_k\}$ and all the values of A are also included. We assume that each of the t individuals belongs to one of k mutually exclusive and exhaustive categories $T = \{T_1, T_2, \dots, T_k\}$, consisting of sensitive categories, $A = \{A_1, A_2, \dots, A_k\}$, and non-sensitive categories, $B = \{B_1, B_2, \dots, B_k\}$, so that $T_i = A_i + B_i$ for $1, 2, \dots, k$. Let t_i denote the random quantity in a category T_i so that $n = \sum_{i=1}^k t_i$. The random quantities a_i and b_i are defined similarly, so that $a = \sum_{i=1}^k a_i$ and $b = \sum_{i=1}^k b_i$, and $t_i = a_i + b_i$. Our goal, then, is to estimate $\pi_1, \pi_2, \dots, \pi_k$, the proportions in the population associated with the sensitive categories, $A = \{A_1, A_2, \dots, A_k\}$. Based on the number of green balls in the device, $p_{b_i} = q_i/g$ is the proportion of green balls with number i for $i = 1, 2, \dots, k$, where q_i is the number of green balls that contain number i , and $g = \sum_{i=1}^k q_i$; that is, the quantities p_{b_i} are known in advance. For a tabular representation of our multinomial situation, see Table 1. Let $p_{t_1}, p_{t_2}, \dots, p_{t_k}$ denote the proportions in the population who are in categories $T = \{T_1, T_2, \dots, T_k\}$. With n different interviewees using the Hopkins' device, b , the total number of people who are in the non-sensitive categories, $B = \{B_1, B_2, \dots, B_k\}$, is a random quantity with expected value $E[b] = ng/(r + g)$, where r denotes the number of red balls in the device. As b_1, b_2, \dots, b_k are also random quantities with expected value $E[b_i] = nq_i/(r + g)$, it follows that $b_k = b - (b_1 + b_2 + \dots + b_{k-1})$. We assume the distributions of T , A , and B are as follows:

$$T = \{T_1, T_2, \dots, T_{k-1}\} \sim \text{Multinomial}(n, p_{t_1}, p_{t_2}, \dots, p_{t_{k-1}}) = \frac{n!}{\prod_{i=1}^k t_i!} \prod_{i=1}^k p_{t_i}^{t_i},$$

$$A = \{A_1, A_2, \dots, A_{k-1}\} \sim \text{Multinomial}(a, \pi_1, \pi_2, \dots, \pi_{k-1}) = \frac{a!}{\prod_{i=1}^k a_i!} \prod_{i=1}^k \pi_i^{a_i},$$

$$B = \{B_1, B_2, \dots, B_{k-1}\} \sim \text{Multinomial}(b, p_{b_1}, p_{b_2}, \dots, p_{b_{k-1}}) = \frac{b!}{\prod_{i=1}^k b_i!} \prod_{i=1}^k p_{b_i}^{b_i},$$

Suppose that $T = A + B$ is fixed and that respondents give truthful answers to both the sensitive and non-sensitive questions. Then, for random quantities a and b , we derive π_i as follows:

$$\pi_i = \frac{(r + g)p_{t_i} - q_i}{r},$$

where $p_{t_i} = t_i/n$.

If a random sample of size n is drawn, and n_i is the number of respondents answering "i", let $\hat{p}_{t_i} = n_i/n$ denote the proportion of respondents answering "i". If $\hat{\pi}_{M_i}$ denotes the estimates of π_i , it follows that

$$\hat{\pi}_{M_i} = \frac{(r + g)\hat{p}_{t_i} - q_i}{r}. \quad (1)$$

3 Bayesian Multinomial Randomized Response Model

Suppose that in each of k categories, individuals are independently classified into one of T_i ($i = 1, \dots, k$) categories. Therefore $T = (T_1, \dots, T_k)$ has a multinomial distribution with parameters n and $\tilde{p}_t = (p_{t_1}, \dots, p_{t_k})$. So it follows that based on the observed values $\tilde{t} = (t_1, \dots, t_k)$, the likelihood function of T given that $\tilde{p}_t = (p_{t_1}, \dots, p_{t_k})$ is

$$f_{T|P_t}(\tilde{t}|\tilde{p}_t) = \frac{n!}{\prod_{i=1}^k t_i} \prod_{i=1}^k p_{t_i}^{t_i},$$

where $p_{t_i} \in \Theta_{r,g,q_i} = (\frac{q_i}{r+g}, \frac{r+q_i}{r+g})$. Based on $p_{t_i} = \frac{r\pi_i + q_i}{r+g}$ in Section 2, we can derive the likelihood function of T given that $\tilde{\pi} = (\pi_1, \dots, \pi_k)$ as follows:

$$f_{T|\Pi}(\tilde{t}|n, r, g, \tilde{q}, \tilde{\pi}) = \frac{n!}{\prod_{i=1}^k t_i} \prod_{i=1}^k \left(\frac{r}{r+g} \pi_i + \frac{q_i}{r+g} \right)^{t_i},$$

where $0 < \pi_i < 1$ and $\tilde{q} = (q_1, \dots, q_k)$.

If $\frac{r}{r+g}$ and $\frac{q_i}{r+g}$ are nonnegative real numbers satisfying $0 \leq \frac{r+q_i}{r+g} \leq 1$, then

$$p_{t_i} = \frac{r\pi_i + q_i}{r+g}, \quad \frac{q_i}{r+g} < p_{t_i} < \frac{r+q_i}{r+g} \text{ and } \sum_{i=1}^k p_{t_i} = 1.$$

and therefore

$$\pi_i = \frac{(r + g)p_{t_i} - q_i}{r}, \quad 0 < \pi_i < 1 \text{ and } \sum_{i=1}^k \pi_i = 1.$$

Thus, the parameter space, which is related to p_{t_i} can be written as

$$\Theta_{r,g,q_i} = \left(\frac{q_i}{r+g}, \frac{r+q_i}{r+g} \right) \subset (0, 1).$$

Hence, proper conjugate prior for $\tilde{p}_t = (p_{t_1}, \dots, p_{t_k})$ is a truncated Dirichlet distribution, denoted by $\text{Dirichlet}_1(\alpha_1, \dots, \alpha_k)$, with density

$$f_{P_t}(\tilde{p}_t | \tilde{\alpha}) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k p_{t_i}^{\alpha_i-1}, \quad (2)$$

where $\alpha_i > 0$ for $i = 1, \dots, k$, $p_{t_i} \in \Theta_{r,g,q_i}$, and Γ is the gamma distribution.. By the linear transformation of $p_{t_i} = \frac{r\pi_i+q_i}{r+g}$, a prior density for $\tilde{\pi} = (\pi_1, \dots, \pi_k)$ is is a modified truncated Dirichlet distribution, denoted by $\text{Dirichlet}_2(\alpha_1, \dots, \alpha_k)$, with density

$$f_{\Pi}(\tilde{\pi} | \tilde{\alpha}, r, g, \tilde{q}) = \left(\frac{r}{r+g} \right)^{k-1} \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k \left(\frac{r}{r+g} \pi_i + \frac{q_i}{r+g} \right)^{\alpha_i-1}, \quad (3)$$

where $\alpha_i > 0$ for $i = 1, \dots, k$, and $\pi_i \in (0, 1)$.

Hence, the posterior distribution of T and P_t is as follows:

$$f_{P_t, T}(\tilde{t}, \tilde{p}_t | \tilde{\alpha}) = \frac{n! \Gamma(\sum_{i=1}^k \alpha_i)}{\left(\prod_{i=1}^k t_i \right) \left(\prod_{i=1}^k \Gamma(\alpha_i) \right)} \prod_{i=1}^k p_{t_i}^{t_i + \alpha_i - 1} \quad (4)$$

for $p_{t_i} \in \Theta_{r,g,q_i}$, and the marginal distribution of T is given by

$$f_T(\tilde{t} | n, \tilde{\alpha}) = \frac{n! \Gamma(\sum_{i=1}^k \alpha_i)}{\left(\prod_{i=1}^k t_i \right) \left(\prod_{i=1}^k \Gamma(\alpha_i) \right)} \frac{\prod_{i=1}^k \Gamma(t_i + \alpha_i)}{\Gamma(\sum_{i=1}^k (t_i + \alpha_i))}. \quad (5)$$

The conditional distribution of P_t given T is

$$f_{P_t | T}(\tilde{p}_t | \tilde{t}, \tilde{\alpha}, n) = \frac{\Gamma(\sum_{i=1}^k (t_i + \alpha_i))}{\prod_{i=1}^k \Gamma(t_i + \alpha_i)} \prod_{i=1}^k p_{t_i}^{t_i + \alpha_i - 1}, \quad (6)$$

which means that $P_t | T$ has a $\text{Dirichlet}_1(t_1 + \alpha_1, \dots, t_k + \alpha_k)$.

Similarly, the posterior distribution of T and Π is as follows:

$$\begin{aligned} f_{\Pi, T}(\tilde{t}, \tilde{\pi} | n, r, g, \tilde{q}, \tilde{\alpha}) &= f_{T | \Pi}(\tilde{t} | n, r, g, \tilde{q}, \tilde{\pi}) \times f_{\Pi}(\tilde{\pi} | \tilde{\alpha}, r, g, \tilde{q}) \\ &= \frac{n! \Gamma(\sum_{i=1}^k \alpha_i)}{\left(\prod_{i=1}^k t_i \right) \left(\prod_{i=1}^k \Gamma(\alpha_i) \right)} \left(\frac{r}{r+g} \right)^{k-1} \prod_{i=1}^k \left(\frac{r}{r+g} \pi_i + \frac{q_i}{r+g} \right)^{t_i + \alpha_i - 1} \end{aligned} \quad (7)$$

for $0 < \pi_i < 1$, and the marginal distribution of T is given by

$$f_T(\tilde{t}|n, r, g, \tilde{q}, \tilde{\alpha}) = \frac{n! \Gamma(\sum_{i=1}^k \alpha_i)}{\left(\prod_{i=1}^k t_i\right) \left(\prod_{i=1}^k \Gamma(\alpha_i)\right) \Gamma(\sum_{i=1}^k t_i + \alpha_i)}. \quad (8)$$

The conditional distribution of Π given T is

$$f_{\Pi|T}(\tilde{\pi}|\tilde{t}, \tilde{\alpha}, n, r, g, \tilde{q}) = \left(\frac{r}{r+g}\right)^{k-1} \frac{\Gamma(\sum_{i=1}^k (t_i + \alpha_i))}{\prod_{i=1}^k \Gamma(t_i + \alpha_i)} \prod_{i=1}^k \left(\frac{r}{r+g} \pi_i + \frac{q_i}{r+g}\right)^{t_i + \alpha_i - 1}, \quad (9)$$

which means that $\Pi|T$ has a Dirichlet $_2(t_1 + \alpha_1, \dots, t_k + \alpha_k)$.

Since $P_t|T$ has a Dirichlet $_1(t_1 + \alpha_1, \dots, t_k + \alpha_k)$, its marginal has a Beta density function as follows:

$$f_{P_{t_i}|T}(p_{t_i}|\tilde{t}, \tilde{\alpha}, n) = \frac{\Gamma(\sum_{i=1}^k (t_i + \alpha_i))}{\Gamma(t_i + \alpha_i) \Gamma(\sum_{j=1}^k (t_j + \alpha_j) - (t_i + \alpha_i))} p_{t_i}^{t_i + \alpha_i - 1} (1 - p_{t_i})^{(\sum_{j=1}^k (t_j + \alpha_j) - (t_i + \alpha_i) - 1)},$$

for $i = 1, \dots, k$. For the remainder of this section, and for all comparisons in Section 5, only squared-error loss is considered; i.e., $L(p_{t_i}, a) = (p_{t_i} - a)^2$, so that the Bayes estimate of p_{t_i} is the mean of the posterior $f_{P_{t_i}|T}(p_{t_i}|\tilde{t}, \tilde{\alpha}, n)$. A simple closed-form expression for $\hat{p}_{t_{B_i}}$, the mean of the posterior, is given by

$$\hat{p}_{t_{B_i}} = E[p_{t_i}|T] = \frac{\text{Beta}(t_i + \alpha_i + 1, \sum_{j=1}^k (t_j + \alpha_j) - (t_i + \alpha_i))}{\text{Beta}(t_i + \alpha_i, \sum_{j=1}^k (t_j + \alpha_j) - (t_i + \alpha_i))}, \quad (10)$$

where

$$\text{Beta}(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}.$$

Since we know $n = \sum_{j=1}^k t_j$ and denote $\beta_i = (\sum_{j=1}^k \alpha_j) - \alpha_i$, (10) reduces to

$$\hat{p}_{t_{B_i}} = \frac{\text{Beta}(t_i + \alpha_i + 1, n - t_i + \beta_i)}{\text{Beta}(t_i + \alpha_i, n - t_i + \beta_i)} \quad (11)$$

The classical estimator, MLE, derived in Section 2, can be obtained from the Bayes estimator (11) by choosing different values of α and β . If $\alpha_1 = \alpha_2 = \dots = \alpha_k = 0$, then the Bayes estimator corresponds to the MLE. Finally, we obtain the estimate of π_{B_i} as follows

$$\hat{\pi}_{B_i} = E[p_{t_i}|T] = \frac{(r+g)\hat{p}_{t_{B_i}} - q_i}{r},$$

and the variance of $\hat{\pi}_{B_i}$ is given by

$$V(\hat{\pi}_{B_i}) = \left(\frac{r+g}{r}\right)^2 \times \left[\frac{\text{Beta}(t_i + \alpha_i + 2, n - t_i + \beta_i)}{\text{Beta}(t_i + \alpha_i, n - t_i + \beta_i)} - \left(\frac{\text{Beta}(t_i + \alpha_i + 1, n - t_i + \beta_i)}{\text{Beta}(t_i + \alpha_i, n - t_i + \beta_i)}\right)^2 \right].$$

4 Application

Using this Bayesian framework, we implement an MCMC method to sample from the full conditional distribution of all these parameters. For squared-error loss and absolute-error loss, WinBUGS provides two Bayes estimates that are the mean and median of the posterior distribution, respectively. We have used four different Dirichlet priors for the sensitive characteristic parameters of the multinomial RR model. Results are based on three chains of 27,000 iterations, each after a burn-in period of 3,000 iterations. Figure 1 denotes the plots of kernel estimates of the marginal posterior density of π_i based on posterior samples using non-information prior, Dirichlet (1,1,1).

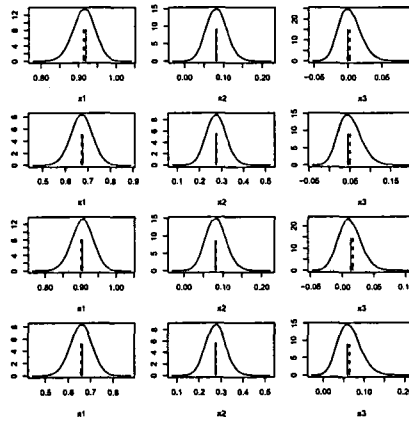


Figure 1: Posterior density of parameters using non-informative Dirichlet prior, Dirichlet (1, 1, 1) (π_{B_i} : Dotted arrow, π_{M_i} : Solid arrow): Case 1 (top row), Case 2 (second row), Case 1* (third row), Case 2* (last row)

References

- [1] Kim, J.-M., Tebbs J., and An, S.-W. (2005), Extensions of Mangat's randomized-response model. *Journal of Statistical Planning and Inference*, In press.
- [2] Kim, J.-M. and Warde, W. D. (2005), Some new results on the multinomial randomized response model. *Communications in Statistics: Theory and Methods*, **34**, 4, xxx-xxx.
- [3] Warner, S. (1965), Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, **60**, 63-69.